

# 確率・統計基礎: KL ダイバージェンスの性質

森 立平

# Kullback–Leibler divergence

## Definition (Kullback–Leibler divergence)

$D: \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}$  を以下で定義する

$$D(q \parallel p) := \sum_{k=1}^m q_k \log \frac{q_k}{p_k}.$$

ただし、 $0 \log 0 = 0 \log \frac{0}{0} = 0$  とし、 $q \log \frac{q}{0}$  は  $q > 0$  について  $+\infty$  とする。

# サノフの定理

## Theorem (サノフの定理)

任意の  $\Gamma \subseteq \mathcal{P}$  について

$$\frac{1}{N} \log \Pr(\hat{P}_{\mathbf{X}} \in \Gamma) \leq - \inf_{q \in \Gamma} D(q \| p) + \frac{m \log(N+1)}{N} \quad \text{for any } N \in \mathbb{Z}_{>0}$$

$$\liminf_{N \rightarrow \infty} \frac{1}{N} \log \Pr(\hat{P}_{\mathbf{X}} \in \Gamma) \geq - \inf_{q \in \Gamma^\circ} D(q \| p).$$

ただし、 $\Gamma^\circ$  は  $\Gamma$  の内点の集合とする。

## Corollary

$\Gamma \subseteq \mathcal{P}$  について  $\Gamma \subseteq \overline{\Gamma^\circ}$  のとき、

$$\lim_{N \rightarrow \infty} \frac{1}{N} \log \Pr(\hat{P}_{\mathbf{X}} \in \Gamma) = - \min_{q \in \overline{\Gamma}} D(q \| p).$$

ただし、 $\overline{\Gamma}$  は  $\Gamma$  の閉包とする。

## KL divergence の性質

1.  $D(q \parallel p) \geq 0$  で等号は  $q = p$  のときのみ成り立つ。
2. 凸関数である。 $D(\lambda q + (1 - \lambda)q' \parallel \lambda p + (1 - \lambda)p') \leq \lambda D(q \parallel p) + (1 - \lambda)D(q' \parallel p')$ .

# イェンセンの不等式

## Lemma (イェンセンの不等式)

凸関数  $f: \mathbb{R} \rightarrow \mathbb{R}$  と確率変数  $X: \Omega \rightarrow \mathbb{R}$  について、

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]).$$

また、 $f$  が狭義凸のとき等号は  $X$  が決定的 ( $\Pr(X = \mathbb{E}[X]) = 1$ ) な場合にのみ成り立つ。

## Proof.

$f(x)$  の  $x = \mathbb{E}[X]$  における“接線”を

$$y = a(x - \mathbb{E}[X]) + f(\mathbb{E}[X])$$

とおくと、 $f$  の凸性より接線は  $f$  より下側にある。つまり

$$f(x) \geq a(x - \mathbb{E}[X]) + f(\mathbb{E}[X]) \quad \forall x \in \mathbb{R}$$

である。 $x = X$  として期待値を取ることによって目的の不等式を得る。

$f$  が狭義凸の場合、 $f$  と接線は一点のみを共有するので決定的でない場合は狭義の不等式が成り立つ。 □

# Log sum 不等式

## Lemma (Log sum 不等式)

任意の  $a_1, \dots, a_m, b_1, \dots, b_m \in \mathbb{R}_{>0}$  について

$$\sum_k a_k \log \frac{a_k}{b_k} \geq \left( \sum_k a_k \right) \log \frac{(\sum_k a_k)}{(\sum_k b_k)}.$$

ある  $c \in \mathbb{R}$  が存在して、任意の  $k \in \{1, \dots, m\}$  について  $a_k = cb_k$  が成り立つ場合のみ等号が成り立つ。

## Proof.

$$\begin{aligned} \sum_k a_k \log \frac{a_k}{b_k} &= \left( \sum_k b_k \right) \sum_k \frac{b_k}{\sum_{k'} b_{k'}} \frac{a_k}{b_k} \log \frac{a_k}{b_k} \\ &\geq \left( \sum_k b_k \right) \left( \sum_k \frac{b_k}{\sum_{k'} b_{k'}} \frac{a_k}{b_k} \right) \log \left( \sum_k \frac{b_k}{\sum_{k'} b_{k'}} \frac{a_k}{b_k} \right) \\ &\quad \text{(} x \log x \text{ の凸性とイェンセンの不等式)} \\ &= \left( \sum_k a_k \right) \log \frac{(\sum_k a_k)}{(\sum_k b_k)}. \end{aligned}$$

□

# KL-divergence の凸性

## Lemma

任意の  $p, p', q, q' \in \mathcal{P}$  と  $\lambda \in [0, 1]$  について

$$D(\lambda q + (1 - \lambda)q' \parallel \lambda p + (1 - \lambda)p') \leq \lambda D(q \parallel p) + (1 - \lambda)D(q' \parallel p').$$

## Proof.

$$\begin{aligned} & (\lambda q_k + (1 - \lambda)q'_k) \log \frac{\lambda q_k + (1 - \lambda)q'_k}{\lambda p_k + (1 - \lambda)p'_k} \\ & \leq \lambda q_k \log \frac{\lambda q_k}{\lambda p_k} + (1 - \lambda)q_k \log \frac{(1 - \lambda)q_k}{(1 - \lambda)p_k} \quad (\text{Log sum 不等式}) \\ & = \lambda q_k \log \frac{q_k}{p_k} + (1 - \lambda)q_k \log \frac{q_k}{p_k}. \end{aligned}$$

両辺を  $k$  について和を取ると目的の不等式を得る。 □

$p' = p$  や  $q' = q$  の場合を考えると、

$$\begin{aligned} D(\lambda q + (1 - \lambda)q' \parallel p) & \leq \lambda D(q \parallel p) + (1 - \lambda)D(q' \parallel p) \\ D(q \parallel \lambda p + (1 - \lambda)p') & \leq \lambda D(q \parallel p) + (1 - \lambda)D(q \parallel p'). \end{aligned}$$

# エントロピー

## Definition (エントロピー)

$H: \mathcal{P} \rightarrow \mathbb{R}$  を以下で定義する

$$H(p) := - \sum_{k=1}^m p_k \log p_k.$$

ただし、 $0 \log 0 = 0$  とする。 $\log$  の底は  $m$  とすることも多い。

$u$  を  $\{1, 2, \dots, m\}$  上の一様分布とすると、 $H(p) = \log m - D(p \parallel u)$  である。



# 多項係数とエントロピー

## Lemma

$N, N_1, \dots, N_m \in \mathbb{Z}_{\geq 0}$  で  $N_1 + \dots + N_m = N$  のとき、

$$\frac{1}{(N+1)^m} e^{NH(p)} \leq \binom{N}{N_1 N_2 \dots, N_m} \leq e^{NH(p)}$$

ただし  $k = 1, 2, \dots, m$  について  $p_k := N_k/N$  とする。

## Proof.

$$\begin{aligned} \frac{1}{(N+1)^m} &\leq \binom{N}{N_1 N_2 \dots, N_m} \prod_{k=1}^m (p_k)^{N p_k} \leq 1 \\ \Leftrightarrow \frac{1}{(N+1)^m} &\leq \binom{N}{N_1 N_2 \dots, N_m} e^{-NH(p)} \leq 1. \end{aligned}$$



# データ圧縮とエントロピー

## データ圧縮の概要

- ▶ データ  $\mathbf{x} \in A^N$  が i.i.d でそれぞれが確率分布  $p \in \mathcal{P}$  に従っていると仮定する。
- ▶ サノフの定理よりデータの経験分布  $\hat{P}_{\mathbf{x}}$  は高い確率で真の分布  $p$  に近い。
- ▶ よってあり得る  $\mathbf{x}$  の数は大体  $\binom{N}{p_1 N p_2 N \dots p_m N} \approx e^{NH(p)}$  である。
- ▶ よって長さ  $\log_m e^{NH(p)} = NH(p)/(\log m)$  の  $A$  の列に圧縮できる。ここで  $H(p)/(\log m)$  はエントロピーの定義の  $\log$  の底を  $m$  にしたものと同じ。

# 確率論振り返り

- ▶ 確率空間  $(\Omega, P)$ . 完全加法性
- ▶ 確率変数、期待値、分散、モーメント
- ▶ 大数の法則、クラメールの定理
- ▶ 中心極限定理
- ▶ サノフの定理

# 課題

▶ 二元エントロピー関数

$h(p) := -p \log_2 p - (1-p) \log_2 (1-p)$  for  $p \in [0, 1]$  のグラフを描け。凸性に気をつけること。